# Sri Akilandeswari women's college, wandiwash

Dr.N.Ambiga
Assistant Professor
Department of Computer Science

# INTRODUCTION TO DATA MINING

## 1: Introduction

---

**Instructor:**
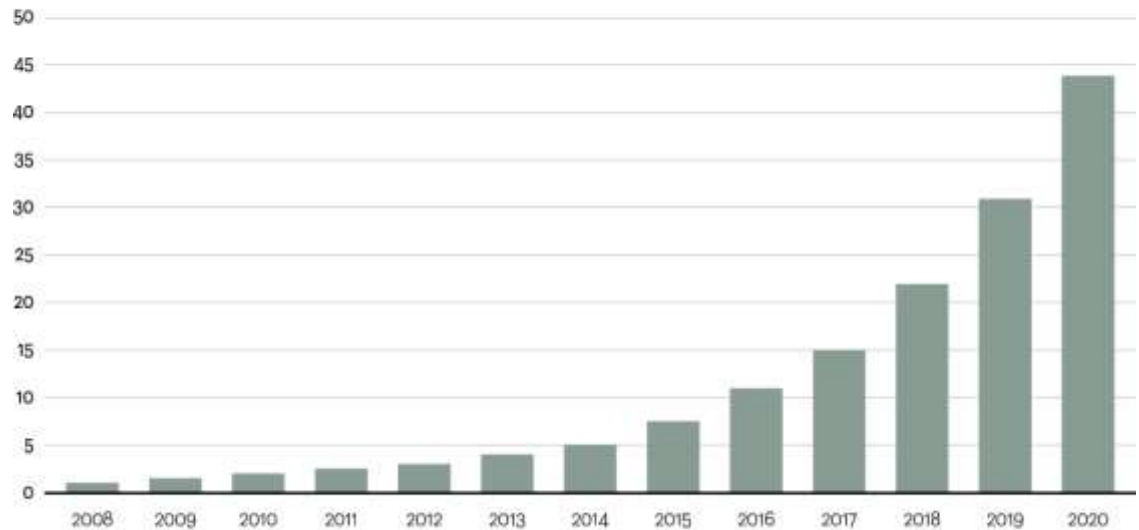
**Dr.N.Ambiga**

# 1. Introduction

- Why Data Mining?

- What Is Data Mining?

- A Multi-Dimensional View of Data Mining

  - What Kinds of Data Can Be Mined?

  - What Kinds of Patterns Can Be Mined?

  - What Kinds of Technologies Are Used?

  - What Kinds of Applications Are Targeted?

- Content covered by this course

# Big Data

- 1 Zeta byte = 1 trillion Gigabytes.

- 5,200 GB of data for every person on Earth.

**Data is growing at a 40 percent compound annual rate, reaching nearly 45 ZB by 2020**

Data in zettabytes (ZB)

| Year | Data (ZB) |
|------|-----------|
| 2008 | 1 |
| 2009 | 1.5 |
| 2010 | 2 |
| 2011 | 2.5 |
| 2012 | 3 |
| 2013 | 4 |
| 2014 | 5 |
| 2015 | 7.5 |
| 2016 | 11 |
| 2017 | 15 |
| 2018 | 22 |
| 2019 | 31 |
| 2020 | 44 |

Source: Oracle, 2012

# Example of Data Volumes

| Unit | Value | Example |
|---|---|---|
| Kilobytes (KB) | 1,000 bytes | a paragraph of a text document |
| Megabytes (MB) | 1,000 Kilobytes | a small novel |
| Gigabytes (GB) | 1,000 Megabytes | Beethoven's 5th Symphony |
| Terabytes (TB) | 1,000 Gigabytes | all the X-rays in a large hospital |
| Petabytes (PB) | 1,000 Terabytes | half the contents of all US academic research libraries |
| Exabytes (EB) | 1,000 Petabytes | about one fifth of the words people have ever spoken |
| Zettabytes (ZB) | 1,000 Exabytes | as much information as there are grains of sand on all the world's beaches |
| Yottabytes (YB) | 1,000 Zettabytes | as much information as there are atoms in 7,000 human bodies |

https://www.eecis.udel.edu/~amer/Table-Kilo-Mega-Giga---YottaBytes.html

# Why Data Mining?

- The Explosive Growth of Data: from terabytes to petabytes

  - Data collection and data availability

    - Automated data collection tools, database systems, Web, computerized society

  - Major sources of abundant data

    - Business: Web, e-commerce, transactions, stocks, …

    - Science: Remote sensing, bioinformatics, scientific simulation, …

    - Society and everyone: news, digital cameras, YouTube, social media, mobile devices, …

- We are drowning in data, but starving for knowledge!

- "Necessity is the mother of invention"—Data mining—Automated analysis of massive data sets

# 1. Introduction

- Why Data Mining?

- What Is Data Mining?

- A Multi-Dimensional View of Data Mining

  - What Kinds of Data Can Be Mined?

  - What Kinds of Patterns Can Be Mined?

  - What Kinds of Technologies Are Used?

  - What Kinds of Applications Are Targeted?
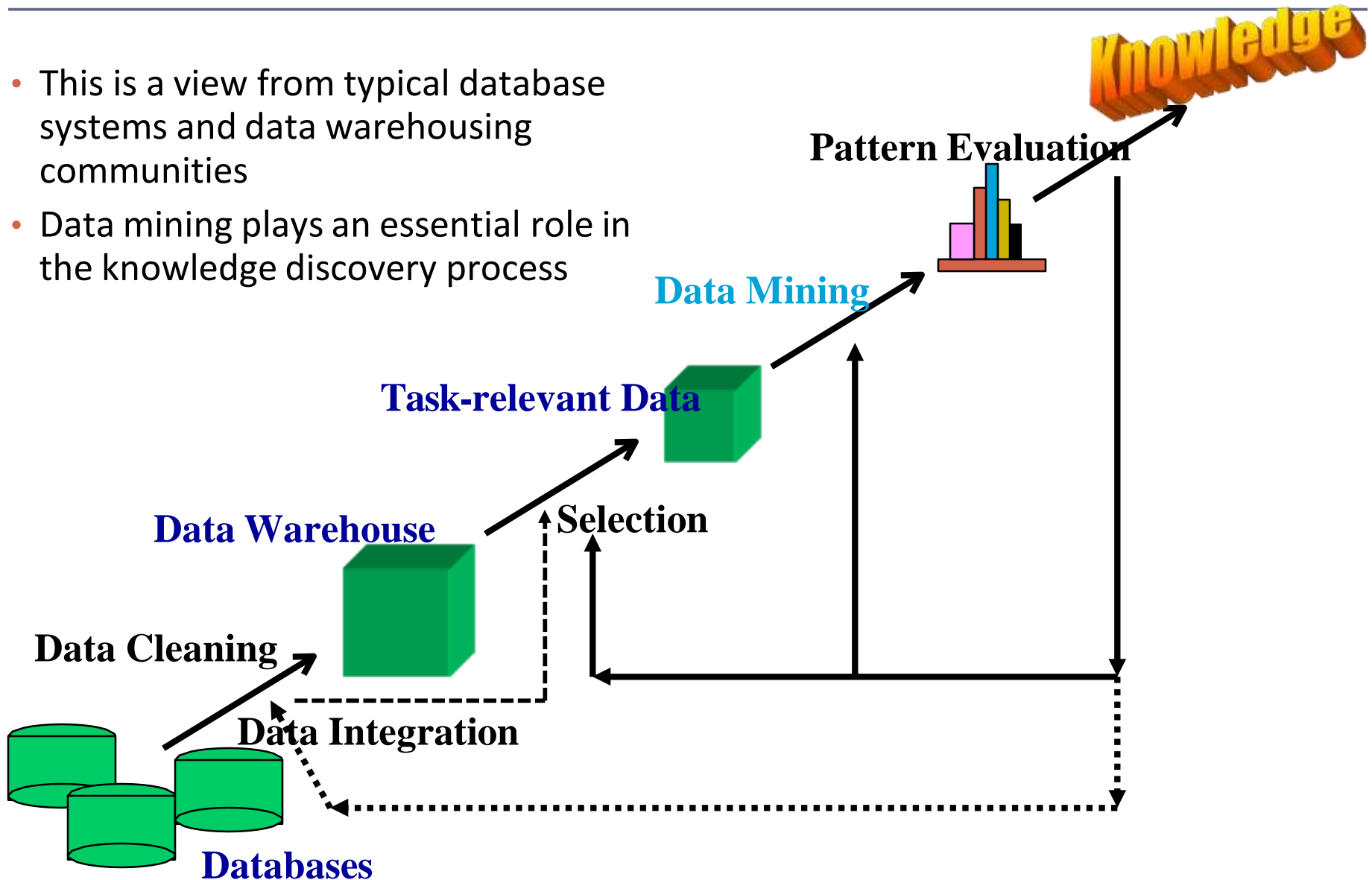
- Content covered by this course

# What Is Data Mining?

- Data mining (knowledge discovery from data)
  - Extraction of interesting (<u>non-trivial,</u> <u>implicit</u>, <u>previously unknown</u> and <u>potentially useful)</u> patterns or knowledge from huge amount of data

- Alternative names
  - Knowledge discovery (mining) in databases (KDD), knowledge extraction, data/pattern analysis, data archeology, data dredging, information harvesting, business intelligence, etc.
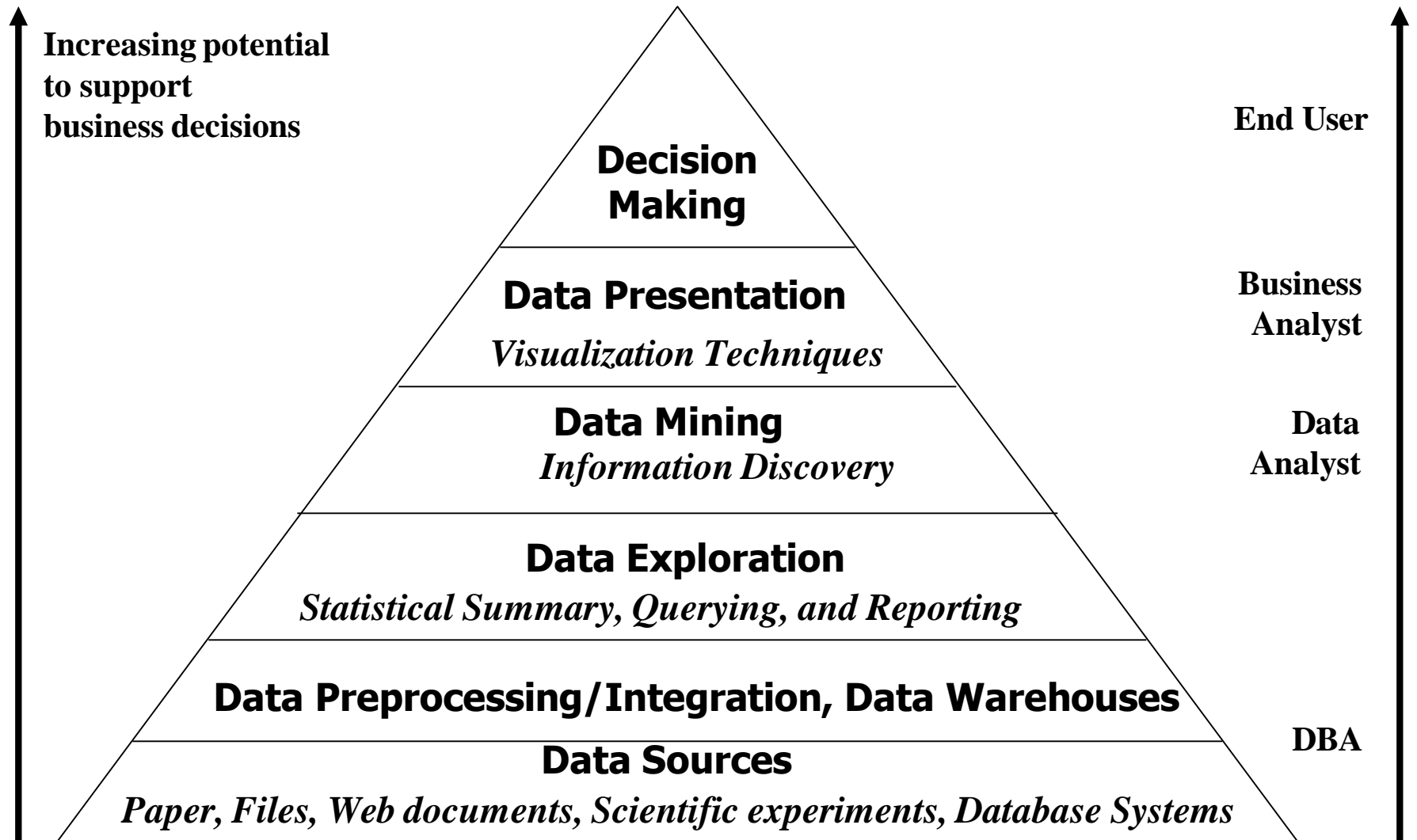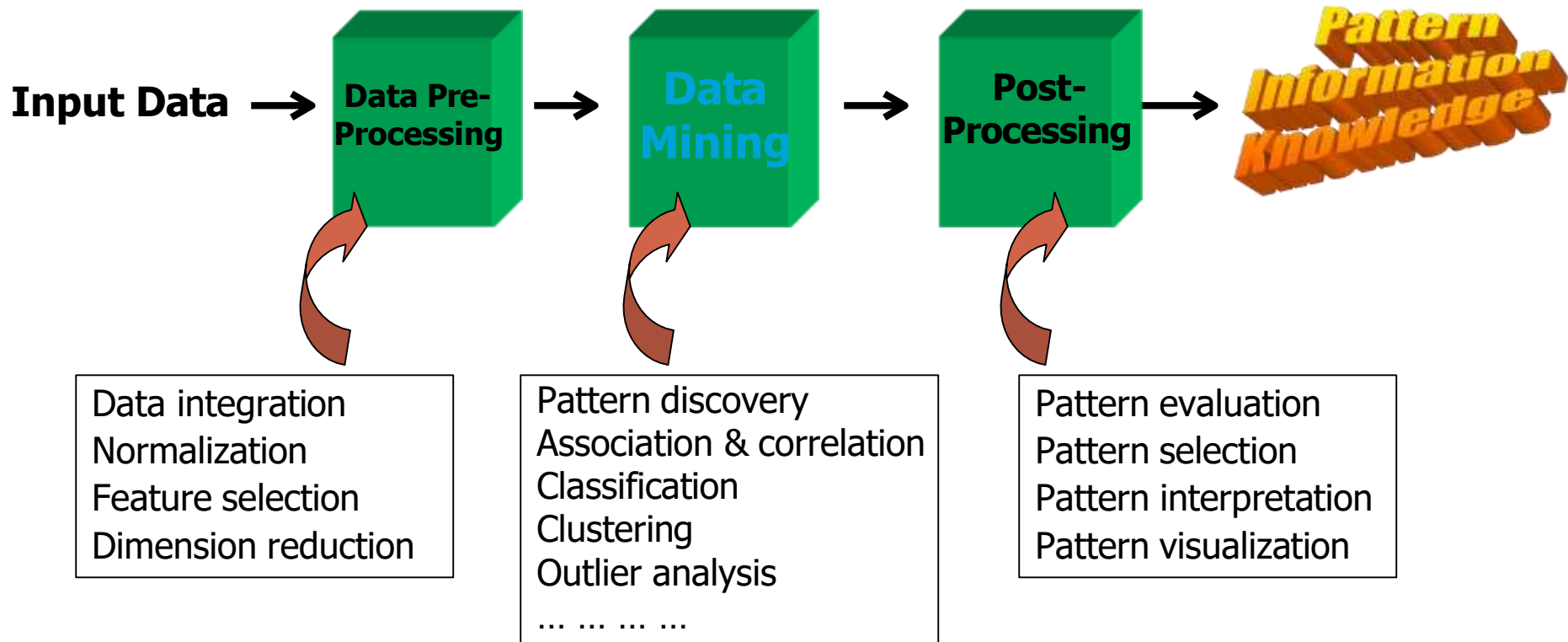
# Knowledge Discovery (KDD) Process

- This is a view from typical database systems and data warehousing communities

- Data mining plays an essential role in the knowledge discovery process

**Knowledge**

**Pattern Evaluation**

**Data Mining**

**Task-relevant Data**

**Data Warehouse**

**Selection**

**Data Cleaning**

**Data Integration**

**Databases**

# Data Mining in Business Intelligence

**Increasing potential
to support
business decisions**

**End User**

**Decision
Making**

**Business
Analyst**

**Data Presentation**

*Visualization Techniques*

**Data
Analyst**

**Data Mining**

*Information Discovery*

**Data Exploration**

*Statistical Summary, Querying, and Reporting*

**Data Preprocessing/Integration, Data Warehouses**

**DBA**

**Data Sources**

*Paper, Files, Web documents, Scientific experiments, Database Systems*

# KDD Process: A Typical View from ML and Statistics

**Input Data** → **Data Pre-Processing** → **Data Mining** → **Post-Processing** → Pattern Information Knowledge

| Data integration | Pattern discovery | Pattern evaluation |
|---|---|---|
| Normalization | Association & correlation | Pattern selection |
| Feature selection | Classification | Pattern interpretation |
| Dimension reduction | Clustering | Pattern visualization |
| | Outlier analysis | |
| | … … … … | |

- This is a view from typical machine learning and statistics communities

# 1. Introduction

- Why Data Mining?

- What Is Data Mining?

- A Multi-Dimensional View of Data Mining

  - What Kinds of Data Can Be Mined?

  - What Kinds of Patterns Can Be Mined?

  - What Kinds of Technologies Are Used?

  - What Kinds of Applications Are Targeted?

- Content covered by this course

# Multi-Dimensional View of Data Mining

- **Data to be mined**
  - Database data (extended-relational, object-oriented, heterogeneous, legacy), data warehouse, transactional data, stream, spatiotemporal, time-series, sequence, text and web, multi-media, graphs & social and information networks
- **Knowledge to be mined (or: Data mining functions)**
  - Characterization, discrimination, association, classification, clustering, trend/deviation, outlier analysis, etc.
  - Descriptive vs. predictive data mining
  - Multiple/integrated functions and mining at multiple levels
- **Techniques utilized**
  - Data-intensive, data warehouse (OLAP), machine learning, statistics, pattern recognition, visualization, high-performance, etc.
- **Applications adapted**
  - Retail, telecommunication, banking, fraud analysis, bio-data mining, stock market analysis, text mining, Web mining, etc.

# 1. Introduction

- Why Data Mining?

- What Is Data Mining?

- A Multi-Dimensional View of Data Mining

  - What Kinds of Data Can Be Mined?

  - What Kinds of Patterns Can Be Mined?

  - What Kinds of Technologies Are Used?

  - What Kinds of Applications Are Targeted?

- Content covered by this course

# Vector/Tabular Data

|       | Sex | Race | Height | Income | Marital Status | Years of Educ. | Liberal-ness |
|-------|-----|------|--------|--------|----------------|----------------|--------------|
| R1001 | M   | 1    | 70     | 50     | 1              | 12             | 1.73         |
| R1002 | M   | 2    | 72     | 100    | 2              | 20             | 4.53         |
| R1003 | F   | 1    | 55     | 250    | 1              | 16             | 2.99         |
| R1004 | M   | 2    | 65     | 20     | 2              | 16             | 1.13         |
| R1005 | F   | 1    | 60     | 10     | 3              | 12             | 3.81         |
| R1006 | M   | 1    | 68     | 30     | 1              | 9              | 4.76         |
| R1007 | F   | 5    | 66     | 25     | 2              | 21             | 2.01         |
| R1008 | F   | 4    | 61     | 43     | 1              | 18             | 1.27         |
| R1009 | M   | 1    | 69     | 67     | 1              | 12             | 3.25         |

# Set Data

| TID | Items |
|-----|-------|
| 1 | Bread, Coke, Milk |
| 2 | Beer, Bread |
| 3 | Beer, Coke, Diaper, Milk |
| 4 | Beer, Bread, Diaper, Milk |
| 5 | Coke, Diaper, Milk |

# Text Data

- "Text mining, also referred to as text data mining, roughly equivalent to text analytics, refers to the process of deriving high-quality information from text. High-quality information is typically derived through the devising of patterns and trends through means such as statistical pattern learning. Text mining usually involves the process of structuring the input text (usually parsing, along with the addition of some derived linguistic features and the removal of others, and subsequent insertion into a database), deriving patterns within the structured data, and finally evaluation and interpretation of the output. 'High quality' in text mining usually refers to some combination of relevance, novelty, and interestingness. Typical text mining tasks include text categorization, text clustering, concept/entity extraction, production of granular taxonomies, sentiment analysis, document summarization, and entity relation modeling (i.e., learning relations between named entities)." –from wiki

# Text Data – Topic Modeling

# Text Data – Word Embedding



king - man + woman = queen

# Sequence Data



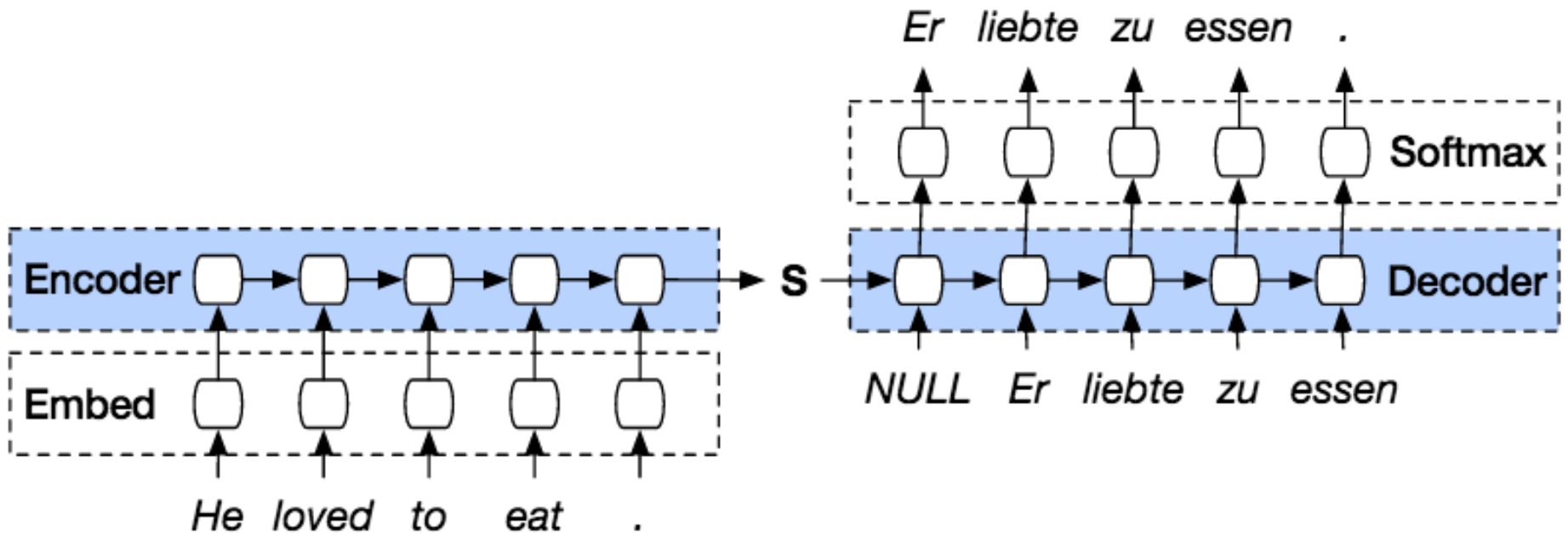SYNTENIC ASSEMBLIES FOR CG15386

```
MD106  ATGCTTAGTAATCCCTACTTTAAGTCCGTTTTGTGGCTGATTGGCTTCGGAGGAATGGG
NEWC   ATGCTTAGTAATCCTTACTTTAAATCCGTTTTGTGGCTGATTGGCTTCGGAGGAATGGG
W501   ATGCTTAGTAATCCCTACTTTAAGTCCGTTTTGTGGCTGATTGGCTTCGGAGGAATGGG
MD199  ATGCTTAGTAATCCCTACTTTAAGTCCGTTTTGTGGCTGATTGGCTTCGGAGGAATGGG
C1674  ATGCTTAGTAATCCCTACTTTAAGTCCGTTTTGTGGCTGATTGGCTTCGGAGGAATGGG
SIM4   ATGCTTAGTAATCCCTACTTTAAGTCCGTTTTGTGGCTGATTGGCTTCGGAGGAATGGG

MD106  CTACGGCCTAATGGTGCTAACAGAGCCGAACGTCGACAAAATAGAGCGCATCAAAGCCT
NEWC   CTACGGCCTAATGGTGCTAACCGAGCCGAACGTCGACAAAATAGAGCGCATCAAAGCCT
W501   CTACGGCCTAATGGTGCTAACCGAGCCGAACGTCGACAAAATAGAGCGCATCAAAGCCT
MD199  CTACGGCCTAATGGTGCTAACCGAGCCGAACGTCGACAAAATAGAGCGCATCAAAGCCT
C1674  CTACGGCCTAATGGTGCTAACCGAGCCGAACGTCGACAAAATAGAGCGCATCAAAGCCT
SIM4   CTACGGCCTAATGGTGCTAACCGAGCCGAACGTCGACAAAATAGAGCGCATCAAAGCCT

MD106  CCGTTTCAAGTACCAAACTGAGTGCGGATGAGCAGCGAAAGGCTCTGTTTATGAAGAAG
NEWC   CCGTTTCAAGTACCAAACTGAGTGCGGATGAGCAGCGAAAGGCTCTGTTTATGAAGAAG
W501   CCGTTTCAAGTACCAAACTGAGTGCGGATGAGCAGCGAAAGGCTCTGTTTATGAAGAAG
MD199  CCGTTTCAAGTACCAAACTGAGTGCGGATGAGCAGCGAAAGGCTCTGTTTATGAAGAAG
C1674  CCGTTTCAAGTACCAAACTGAGTGCGGATGAGCAGCGAAAGGCTCTGTTTATGAAGAAG
SIM4   CCGTTTCAAGTACCAAACTGAGTGCGGATGAGCAGCGAAAGGCTCTGTTTATGAAGAAG

MD106  CTGCAGGAGGCGTCCACCACCAGTGCCCCAATCTACAGGTCAGCGGCCGAGAAATAG
NEWC   CTGCAGGAGGCGTCCACCACCAGTGCCCCAATCTACAGGTCATCGGCCGAGAAATAG
W501   CTGCAGGAGGCGTCCACCACCACTGCCCCAATCTACAGGTCATCGGCCGAGAAATAG
MD199  CTGCAGGAGGCGTCCACCACCAGTGCCCCAATCTACAGGTCAGCGGCCGAGAAATAG
C1674  CTGCAGGAGGCGTCCACCACCAGTGCCCCAATCTACAGGTCAGCGGCCGAGAAATAG
SIM4   CTGCAGGAGGCGTCCACCACCAGTGCCCCAATCTACAGGTCAGCGGCCGAGAAATAG
```
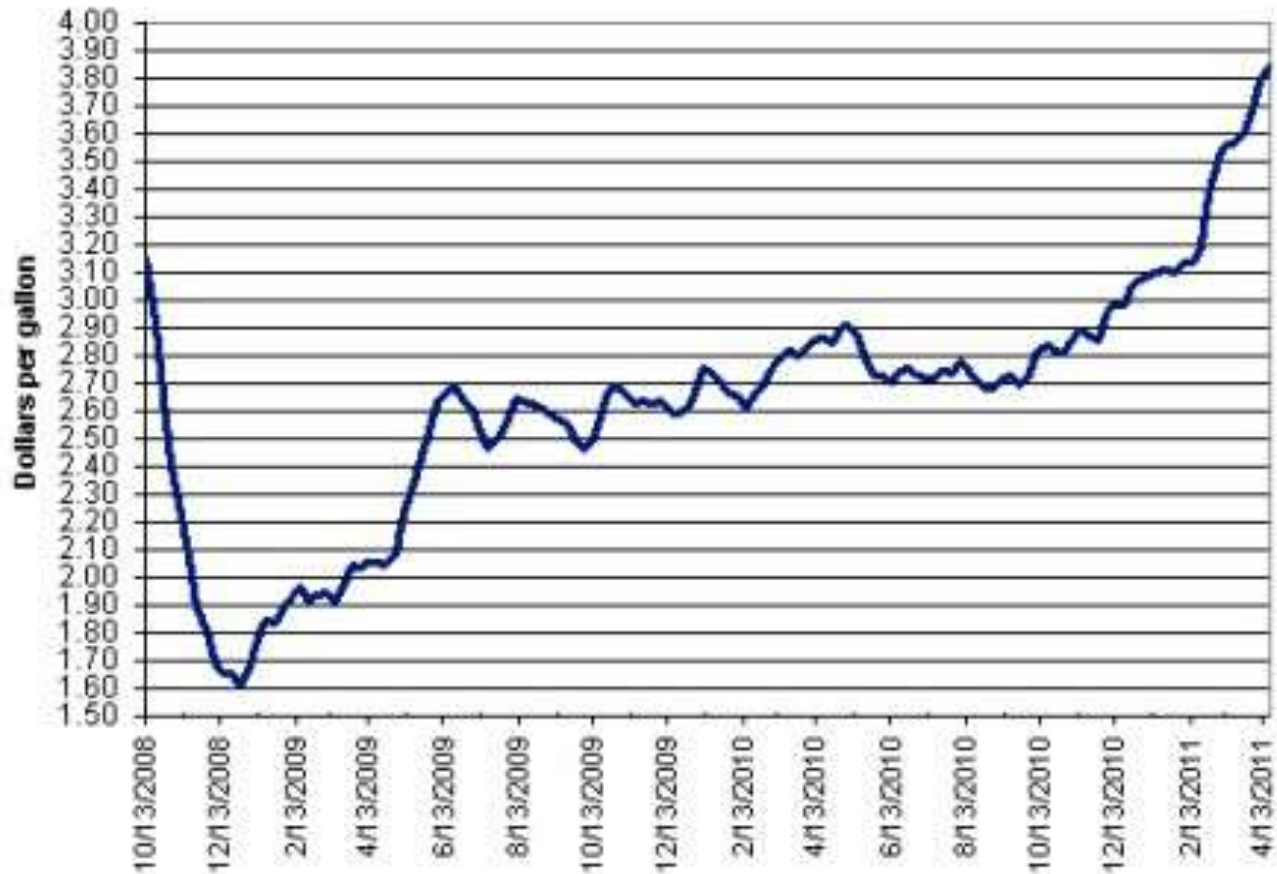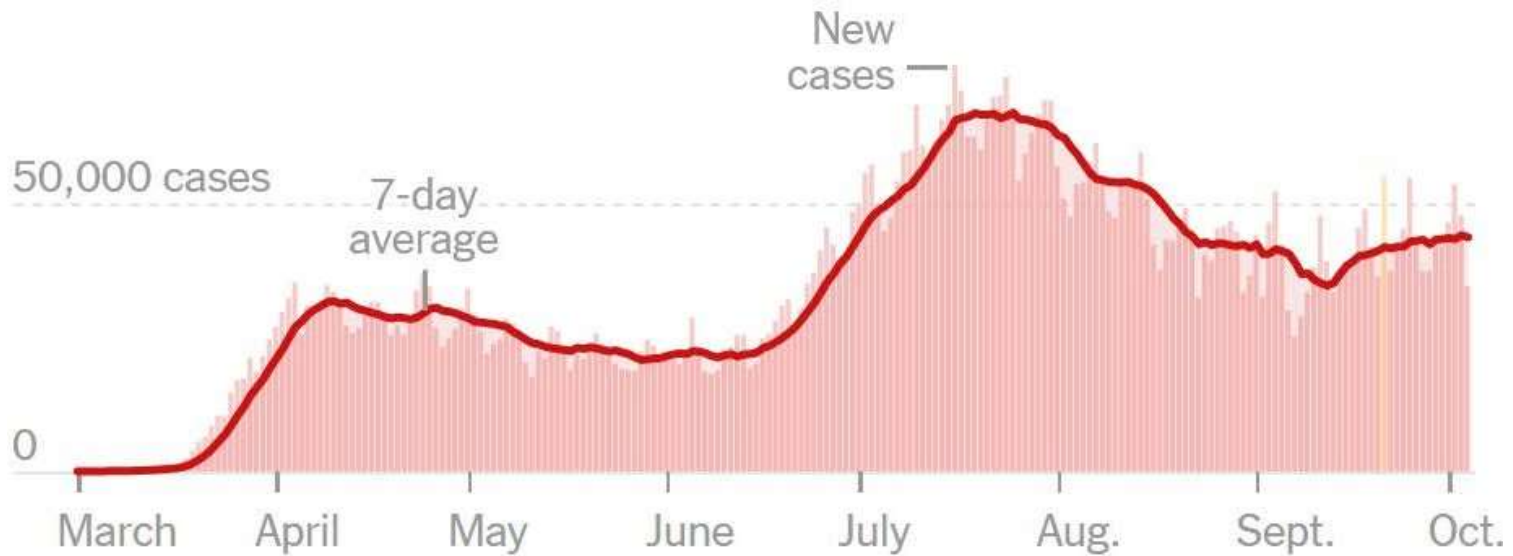
# Sequence Data – Seq2Seq

# Time Series



Weekly U.S. Retail Gasoline Prices, Regular Grade
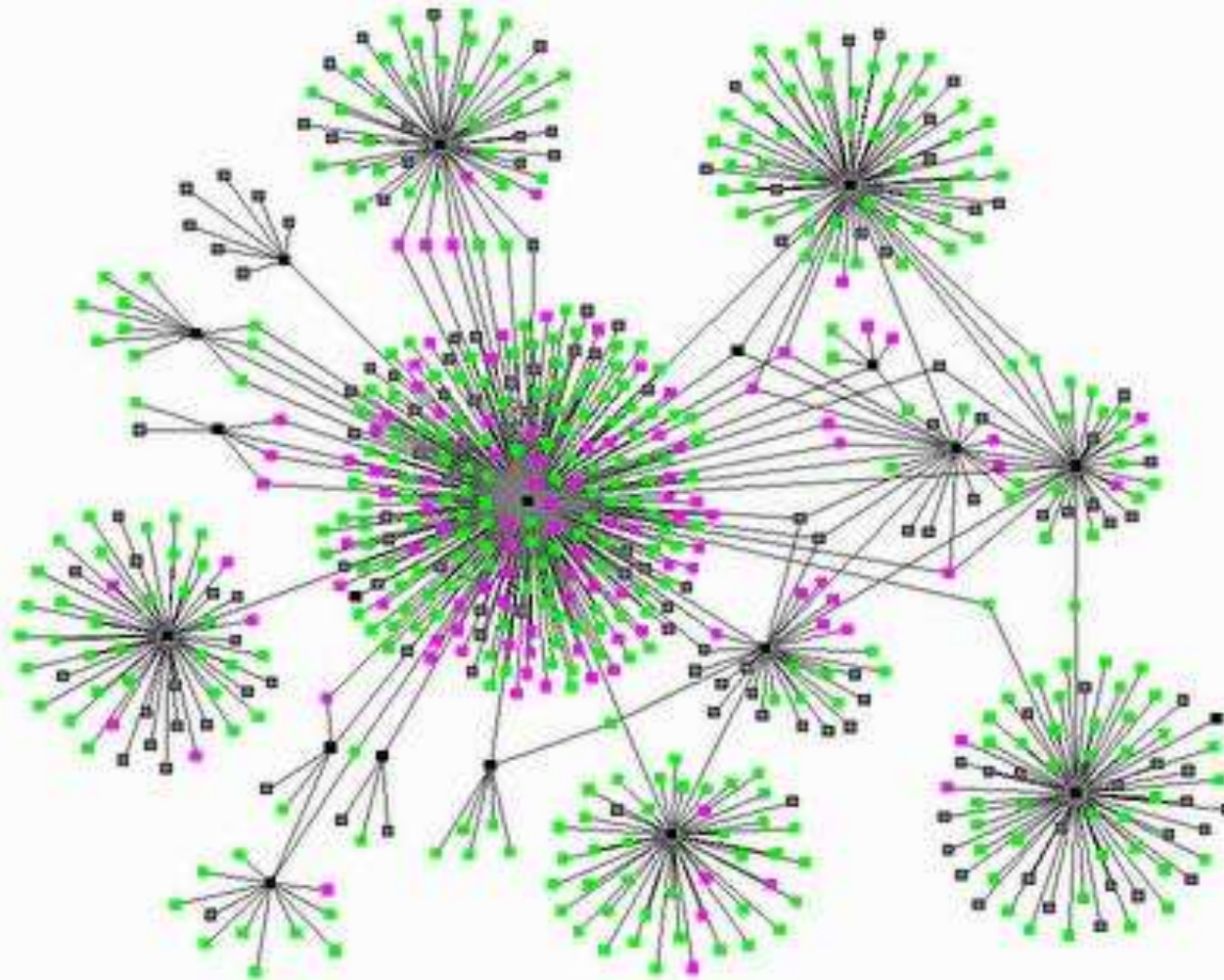
Source: Energy Information Administration

New cases

50,000 cases
7-day average

0

March    April    May    June    July    Aug.    Sept.    Oct.

| | TOTAL REPORTED | ON OCT. 4 | 14-DAY CHANGE |
|---|---|---|---|
| **Cases** | **7.4 million+** | **34,491** | **+6%** → |
| **Deaths** | 209,603 | 332 | −8% → |

■ Day with data reporting anomaly.

Includes confirmed and probable cases where available. 14-day change trends use 7-day averages.
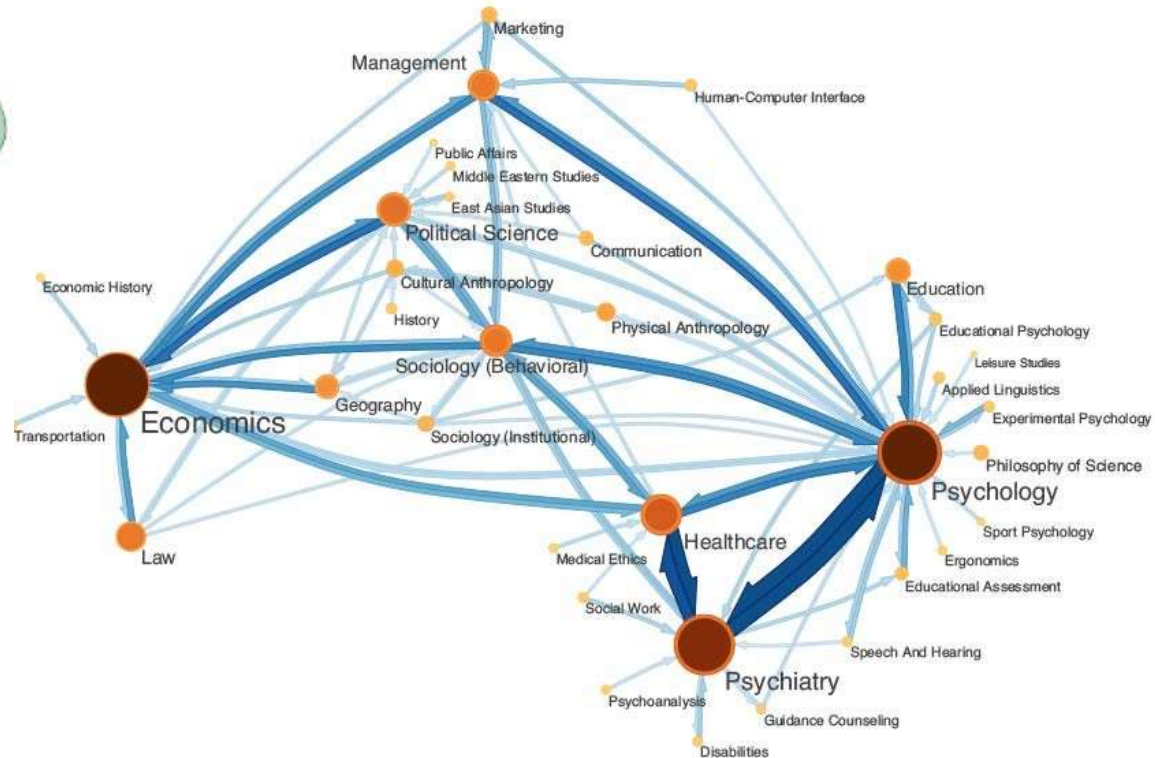
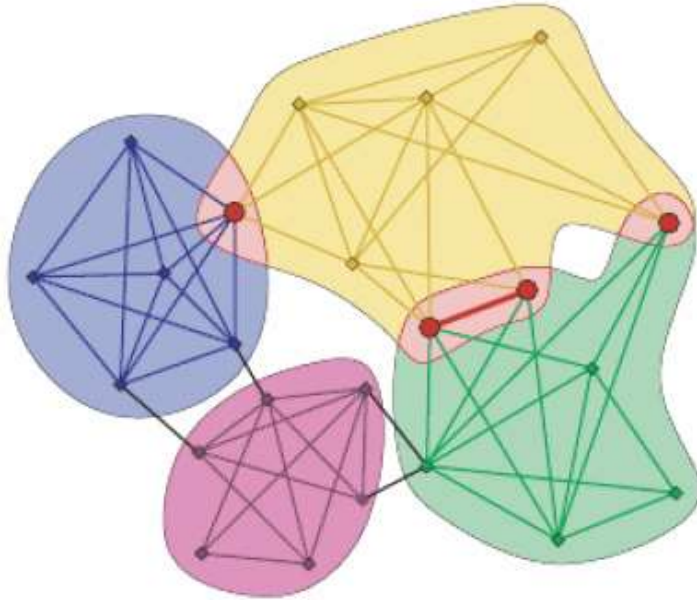# Graph / Network

# Graph / Network – Community Detection

# Image Data

# Image Data – Neural Style Transfer

# Image Data – Image Captioning



"man in black shirt is playing guitar."

"construction worker in orange safety vest is working on road."

"two young girls are playing with lego toy."

"girl in pink dress is jumping in air."

"black and white dog jumps over bar."

"young girl in pink shirt is swinging on swing."

# 1. Introduction

- Why Data Mining?

- What Is Data Mining?

- A Multi-Dimensional View of Data Mining

    - What Kinds of Data Can Be Mined?

    - What Kinds of Patterns Can Be Mined?

    - What Kinds of Technologies Are Used?

    - What Kinds of Applications Are Targeted?

- Content covered by this course

- Frequent patterns (or frequent itemsets)

  - ## What items are frequently purchased together in your Amazon transactions?

Frequently bought together

Total price: $105.88

Add all three to Cart

Add all three to List

- Association, correlation vs. causality

  - ## A typical association rule

    - Diaper $\rightarrow$ Beer [0.5%, 75%]  (support, confidence)

# Data Mining Function: Classification

- Classification and label prediction

  - Construct models (functions) based on some training examples

  - Describe and distinguish classes or concepts for future prediction

    - E.g., classify countries based on (climate), or classify cars based on (gas mileage)

  - Predict some unknown class labels

- Typical methods

  - Decision trees, naïve Bayesian classification, support vector machines, neural networks, rule-based classification, pattern-based classification, logistic regression, …

- Typical applications:

  - Credit card fraud detection, direct marketing, classifying stars, diseases,  web-pages, …
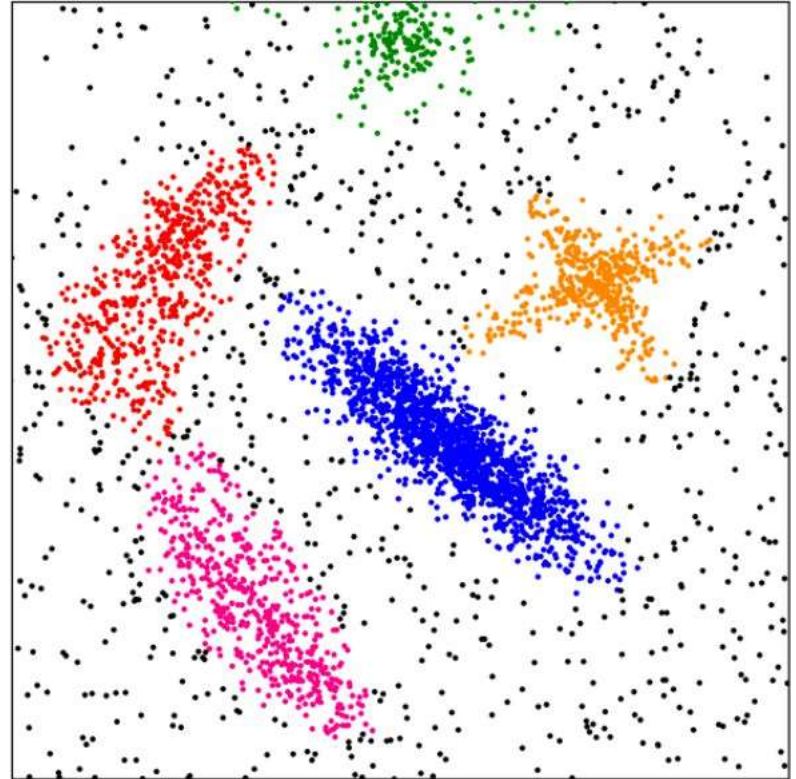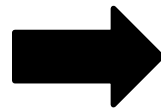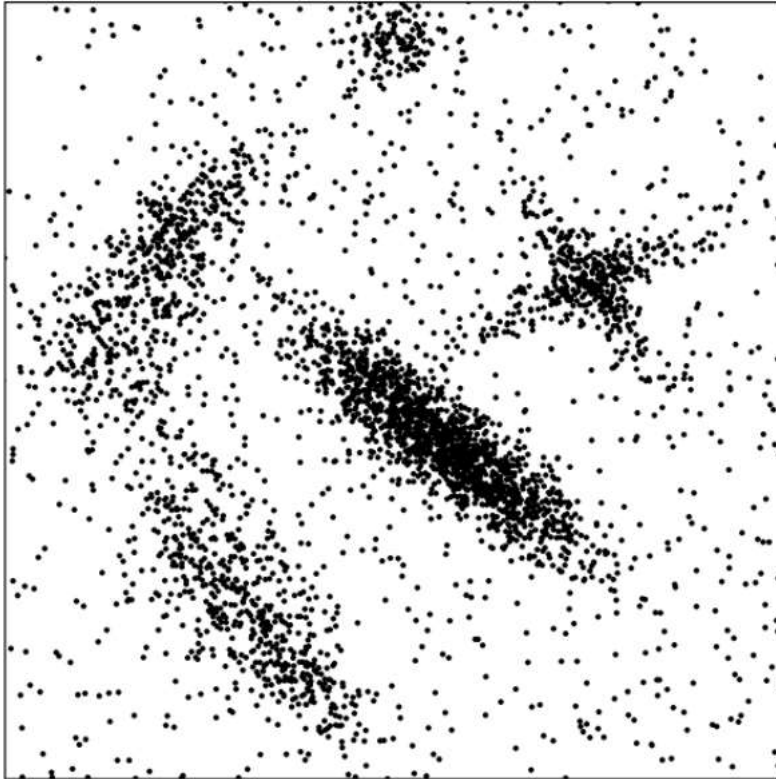
# Image Classification Example

# Data Mining Function: Cluster Analysis

- Unsupervised learning (i.e., Class label is unknown)
- Group data to form new categories (i.e., clusters), e.g., cluster houses to find distribution patterns
- Principle: Maximizing intra-class similarity & minimizing interclass similarity
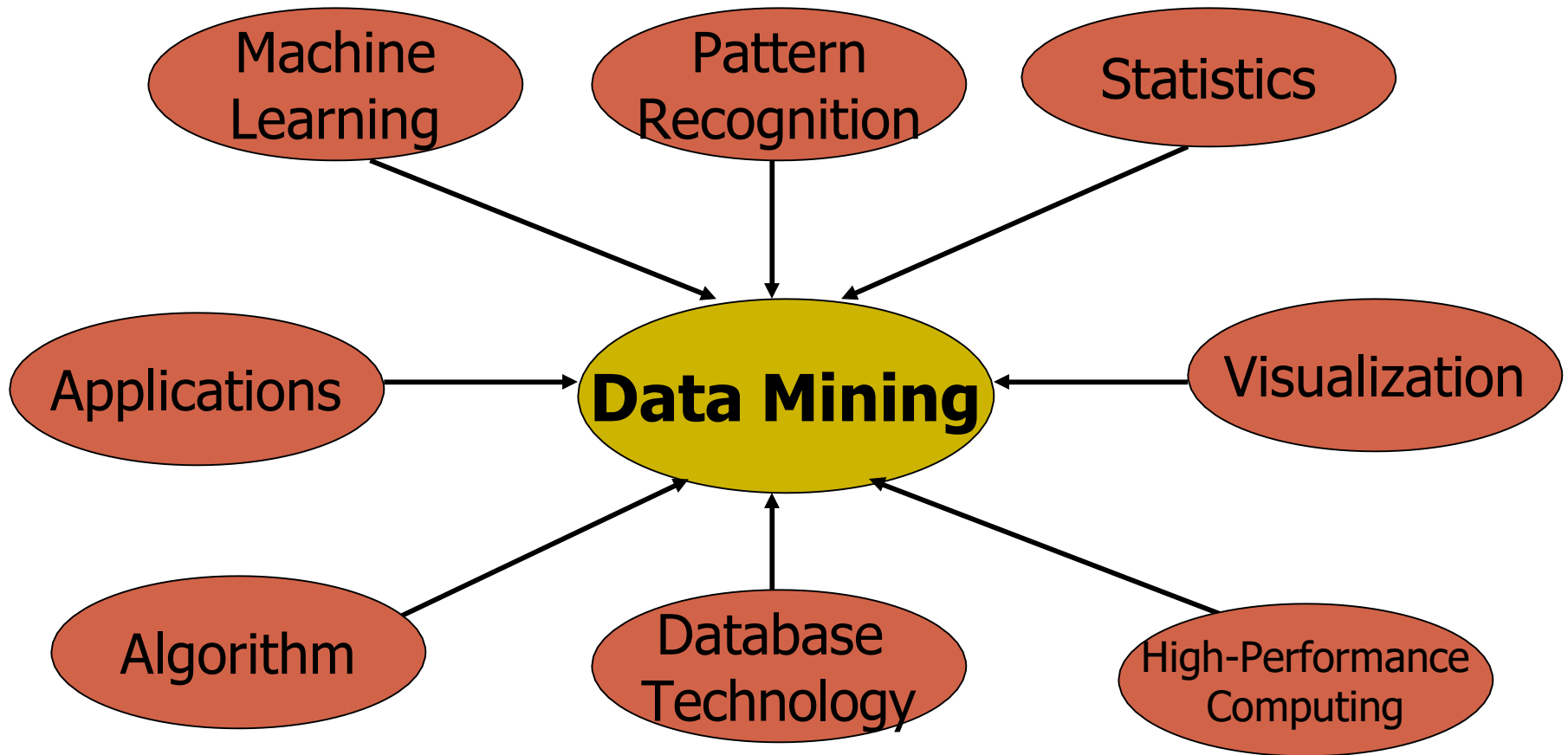- Many methods and applications

# Clustering Example

# Data Mining Functions: Others

- Prediction

- Similarity search

- Ranking

- Outlier detection

- …

# Data Mining: Confluence of Multiple Disciplines

# Applications of Data Mining

- Web page analysis: from web page classification, clustering to PageRank & HITS algorithms

- Collaborative analysis & recommender systems

- Basket data analysis to targeted marketing

- Biological and medical data analysis: classification, cluster analysis (microarray data analysis),  biological sequence analysis, biological network analysis

- Data mining and software engineering (e.g., IEEE Computer, Aug. 2009 issue)

- Social media

- Game

# Thank you!!!!